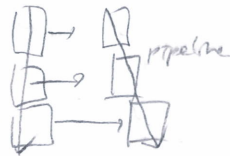


Prog. Model 1: Sequential (SISD)

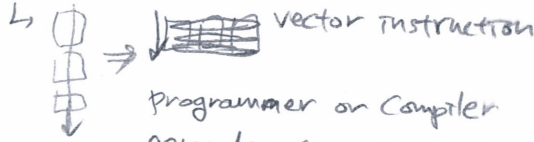


↳ pipelined processor

Out of order execution processor

(← the loop is dynamically unrolled)  
by the hardware

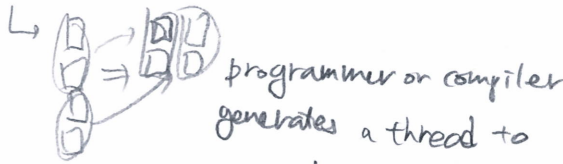
Prog. Model 2: Data Parallel (SIMD)



programmer or compiler  
generates SIMD instruction!

→ Best executed by SIMD processor

Prog. Model 3: Multi-threaded



programmer or compiler  
generates a thread to  
execute each iteration.

→ By MIMD machine.

→ AKA: SPMD

(single program multiple data)

A GPU is a SIMD (SIMT) machine

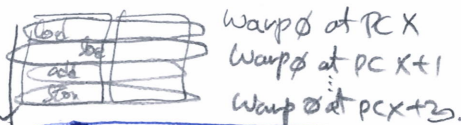
↳ not programmed using "SIMD instr"

• Programmed using threads (SPMD)

• A set of thread executing same instruction  
are dynamically grouped into a Warp  
(wavefront) by the hardware.

→ Warp = SIMD operation by the hardware  
(set of threads executing same instr)

• SIMT: Single instruction multiple threads.



SPMD on SIMT machine

• SIMD: A single sequential INS stream of SIMD  
Instr.  
[VLD, VLD, VADD, VST], VLEN

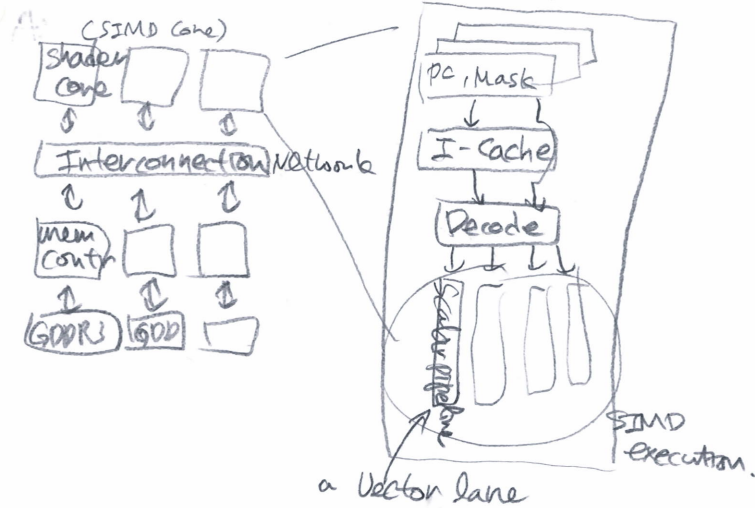
• SIMT: multiple INS streams of scalar instr  
→ grouped to warp  
[LD, LD, ADD, ST], NumThread.

Assume a warp = 32 threads

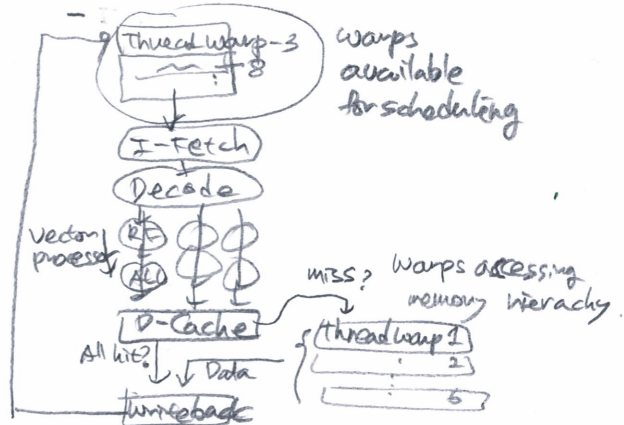
if  $N = 32K \rightarrow 1000$  warps  
(1 iter/thread)

• Warps can be interleaved on the same  
pipeline → Fine grained multithreading  
of warps.

< High level view of GPU >



Latency hiding via Warp-level FGMT.



## Lect. 22 GPU Programming

- Nvidia GeForce GTX 285C (2009)

Nvidia-speak: 240 stream processors  
"SIMT" execution

- ↳ Generic - speak: 30 Cores,  
8 SIMD functional units per core  
(=8 vector lanes)

32 thread  $\rightarrow$  warp  $\rightarrow$  32 warps in FGMT manner

- Nvidia V100  $\rightarrow$  560 stream processor manner.  
 $= 80$  cores w/ 64 SIMD func.  
 $\oplus$   $\dots$  units

① Specialized Functional Units  
for machine learning  
("tensor"-cores)

15.1 TFLOPS : single precision

1.2 Tflops : double prec.

125 TFLOPS for Deep Learning ("Terror" cores)



- Inherent Parallelism → Matrices  
Deep learning  
Image processing  
New programming tool

- Bottleneck

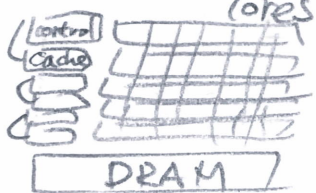
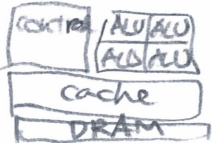
→ CPU-GPU data transfer (PCIe, NVLINK)

- ↳ DRAM memory bandwidth (GDDR5, HBM2)
- ↳ Data layout.

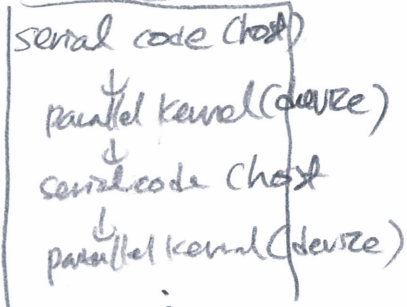
- CPU VS GPU

↳ many in-order FSM T

A. few out of order cores



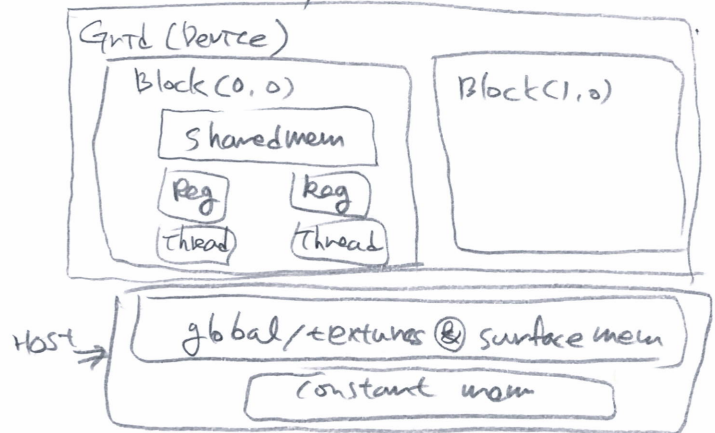
- GPU Computing / offload



- **CUDA / OpenCL** programming model.

- Bulk synchronous programming
  - ↳ Global (Coarse-grain) sync.
- Host (typically CPU) allocates memory, copies data, launch kernels.
- Device (GPU) executes kernels
  - ↳ Grid (NDRange) → Block (work-group)  
→ Thread (work-item)

- Memory Hierarchy





## • CUDA programming

- Fun prototypes: floats serialFunction(...);

--global-- void kernel(...);

- main()

1) Allocate memory  $\rightarrow$  cudaMalloc

2) Data from Host to Device  $\rightarrow$  cudaMemcpy

3) Set #blocks & #threads

4) Kernel Call  $\Rightarrow$  kernel<<execution configuration>>(args);

5) Result from Device to Host  $\rightarrow$  cudaMemcpy

- kernel  $\Rightarrow$  --global-- void kernel ( )

$\hookrightarrow$  registers

shared memory: --shared--

Intra-block sync: \_\_syncthreads()

- Mem deallocation  $\Rightarrow$  cudaFree(d-in);

- Explicit sync  $\rightarrow$  cudaDeviceSynchronize();

• image layout = ex Row-major layout in (1D).

$\hookrightarrow$  one GPU thread per pixel

Grid of Blocks of Threads

$\Rightarrow$  gridDim.x, blockDim.x

$\oplus$  blockIdx.x, threadIdx.x

$\hookrightarrow$  pixel address =  $\text{blockIdx.x} \times \text{blockDim.x} + \text{threadIdx.x}$

• (2D) Grid  $\Rightarrow$  gridDim.x, gridDim.y  
blockDim.x, blockDim.y  
blockIdx.x, blockIdx.y

## [ Memory Access ]

① Latency Hiding w/ FGMT.

② Occupancy

③ Memory Coalescing

AoS vs SoA

(Array of Structure) (Structure of Array)

④ Data reuse: tiling  $\rightarrow$  shared memory.

$\hookrightarrow$  --shared-- int l\_data[(L-size+1)\*(L-size+1)];

⑤ Shared memory = banked (interleaved) mem.

$\hookrightarrow$  32 banks; Bank = Address % 32

$\hookrightarrow$  Bank Conflict  $\rightarrow$  Padding

Randomized mapping

Hash fun.

## [ SIMD Utilization ]

• Intra warp divergence

$\hookrightarrow$  ex Vector reduction: Naïve Mapping

$\downarrow$   
Divergence free mapping

## [ Atomic operations ]

• Atomic conflicts

Histogram Calculation..

• Privatization: Per-block sub-histograms in shared memory

## [ Data Transfers ]

synchronous

Asynchronous

$\hookrightarrow$  Divide into N streams.

$\hookrightarrow$  Video processing

INDG:

**Goal**: Design an acceleration that has

- Simple / Regular Design / High Concurrency / Balanced Computation & I/O (mem) bandwidth

Idea: Replace a single PE with a regular array of PEs and carefully orchestrate flow of data between PEs

Schematischer:  instead of this



instead of this



Sytoliz

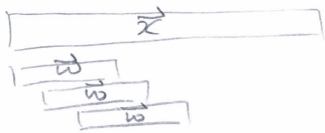
MISD

(memory  $\rightarrow$  heart.  
Data  $\rightarrow$  blood  
pEs  $\rightarrow$  cells)

- Differences in pipelining  $\rightarrow$  Array of PEs

- (Can be non-linear & multi-dimensional)

ex) Convolution.



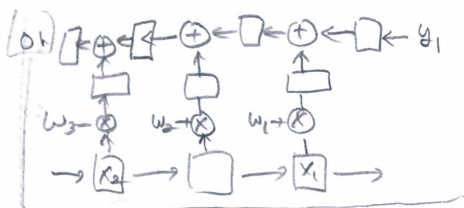
$$\begin{aligned} y_1 &= w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots \\ y_2 &= w_2 x_2 + w_3 x_3 + w_4 x_4 + \dots \\ &\vdots \end{aligned}$$

(-)

$y_{out} = y_{in} + w \cdot x_{in}$   
 $x_{out} = x_{in}$

The diagram illustrates a single neuron and a sequence of neurons. The single neuron is represented by a square box labeled 'W'. It has two inputs:  $x_{in}$  (bottom-left) and  $y_{in}$  (top-right). It has two outputs:  $x_{out}$  (bottom-right) and  $y_{out}$  (top-left). Below this, a sequence of four neurons is shown, each in a square box labeled  $w_1, w_2, w_3, w_4$  from left to right. The first neuron  $w_1$  has an input  $x_3$  and an output  $y_1$ . The second neuron  $w_2$  has an input  $x_2$  and an output  $y_2$ . The third neuron  $w_3$  has an input  $x_1$  and an output  $y_3$ . The fourth neuron  $w_4$  has an input  $y_1$  and an output  $y_4$ . A curved arrow labeled '2 cycle' connects the output  $y_1$  of the first neuron to the input  $y_1$  of the fourth neuron, indicating a two-cycle delay.

$w_i$ 's stay <sup>fixed</sup> &  $x_i$  &  $y_i$  moves systemically  
in opposite direction.



- Carefully Orchestrate !
- Two Dimensional Sparse Arrays & Combination of two different Sparse Array.

Advantages  $\rightarrow$  principle: w/ limited memory bandwidth  
balances comp. to I/O

Specialize: efficiency, simple, high concurrency.

Downside  $\rightarrow$  specialize

more generality  $\Rightarrow$  multiple weights in a RE

+ Data memory in PE  
(temp. result)

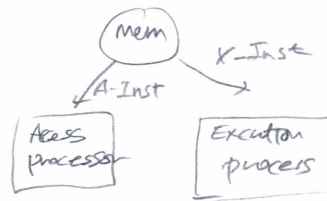
→ pipeline-parallel Program

put stages in a loop  
in a pipeline.

- The Warp Computer.

- Modern Systolic Array : TPU

### Decoupled Access/Execute (DAE)



# Lec 23 b. Memory Organization & Technology



- System maps virtual mem. addresses to physical mem.

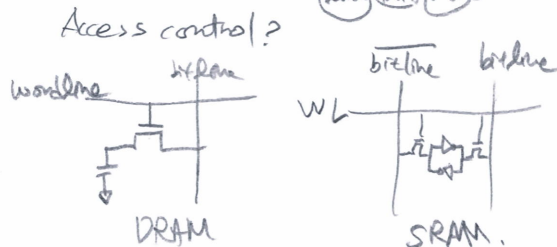
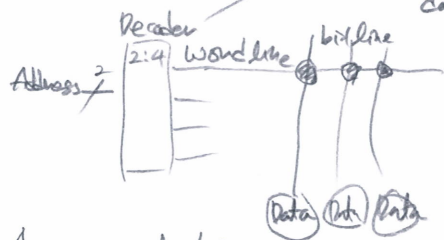
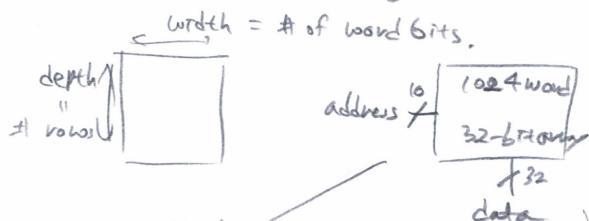


## Devices

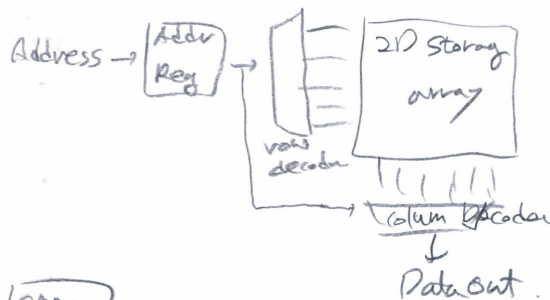
- Flip-Flops (or Latches) ; fast, expensive, parallel > 10 Tr
- Static-RAM ; fast, expensive, one data word = 6 Tr at a time
- Dynamic-RAM ; slower, one data at a time. refresh, deep 1 Tr + 1 Cap
- Storage (Flash, HDD) ; very cheap, non-volatile

Array = memory array

- + Address selection logic
- + Read out circuitry



- Make larger? → Divide the mem into smaller arrays.
- DRAM ⇒ Channel → Rank → Bank  
↓  
Sub arrays → Interleaving  
↓  
Mats



## SRAM

- Address decode
  - drive row select
  - selected bit-cells drive bitlines
  - entire row is read together
  - Differential Sensing & Column Select
  - precharge all bitlines
- Latency

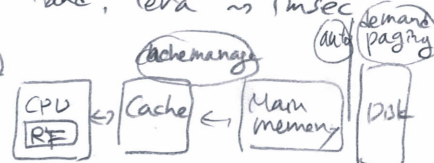
## DRAM

- Read →
- Address decode
  - drive row select
  - selected bit-cell drive bitlines
  - A "flip-flopping" sense amplifier & regen the bitline
  - precharge all bitlines.
- ⇒ Destructive read

## Speed

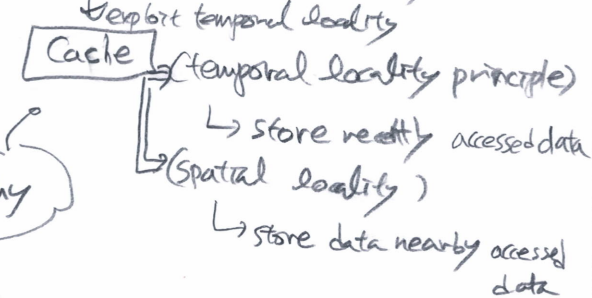
- SRAM, 512 Byte → subnano sec < 10 \$/MB
- SRAM, KmByte → A few Sec
- DRAM, Gigabyte → 50 nsec < 1 \$/MB
- Hard, Tera ~ 1 msec < 1 \$/GB

## Hierarchy



## Memory Locality

in temporal, spatial way



## Cache

- Block, HIT vs Miss, Placement, Replacement
- Granularity of management, Write Policy, Instruction/data bundle.
- Direct-Mapped Cache & Associativity
- LRU
- Shared Cache

